

Multi-view learning for multivariate performance measures optimization

Jim Jing-Yan Wang

Abstract—In this paper, we propose the problem of optimizing multivariate performance measures from multi-view data, and an effective method to solve it. This problem has two features: the data points are presented by multiple views, and the target of learning is to optimize complex multivariate performance measures. We propose to learn a linear discriminant functions for each view, and combine them to construct a overall multivariate mapping function for multi-view data. To learn the parameters of the linear discriminant functions of different views to optimize multivariate performance measures, we formulate a optimization problem. In this problem, we propose to minimize the complexity of the linear discriminant functions of each view, encourage the consistences of the responses of different views over the same data points, and minimize the upper boundary of a given multivariate performance measure. To optimize this problem, we employ the cutting-plane method in an iterative algorithm. In each iteration, we update a set of constrains, and optimize the mapping function parameter of each view one by one.

Index Terms—Multivariate performance measures, Multi-view learning, Cutting-plane algorithm

I. INTRODUCTION

DIFFERENT multivariate performance measures are used to evaluate different machine learning applications [1]. For example, in problems of text classification, F1-score and precision/recall breakeven point (PRBEP) are used to compare true class labels against predicted class labels of a given test text set. In image retrieval problems, the precision/recall at the top k returned images are used to evaluate the performance of a retrieval system. Recently, the problem of multivariate performance measures optimization are proposed to learn directly to losses based on pre-defined multivariate performance measures. Many algorithms have been proposed to solve this problem, and some examples are listed as follows.

- Joachims [1] proposed a support vector machine (SVM)-based for multivariate nonlinear performance measures optimization problems. The proposed algorithm can train a multivariate SVM in polynomial time for large classes of potentially non-linear performance measures. Moreover, the traditional SAM can be arised as a special case of this method.
- Li et al. [2] proposed a two-step approach to optimize multivariate performance measures, by first training a nonlinear classifiers with existing learning methods, and then adapting it to optimize specific performance measures. In the seconde step, the classifier adaptation problem can be solved as a quadratic program problem, in a similar way to linear SVM.
- Mao and Tsang [3] proposed a novel feature selection method to optimize multivariate performance measures,

by formulating the problem or high-dimensional data, and employing a two-layer cutting plane algorithm to solve it. Moreover, this method is also used in multiple-instance learning problems.

Up to now, all these multivariate methods are limited to learn from data with single view. However, in many real-world machine learning problems, the data can be presented by multiple views. For example, in computer vision problems, we can extract different types of features, such as color features, texture features, and shape features, and each feature can be treated as a view. In scientific article classification problems, we can also learn from the views of article abstract, content, and references. Different views of data may be complementary to each other in a learning problem and using multiple views has been a popular strategy in machine learning community. To learn from multiple views of data, different multi-view learning methods have been proposed [4], [5], [6], [7]. However, none of them are designed to optimized a specific multivariate performance measure.

To overcome this problem, in this paper, we propose the problem of learning from multiple view data to optimize the multivariate performance measures. Given a tuple of data points, each of them are presented with multiple views. The problem is to learn a multivariate mapping function to map them to a tuple of class labels, so that the multivariate performance measures can be optimized. To solve this problem, we proposed a novel method for multi-view learning to optimize multivariate performance measures. The contribution of this paper are of two folds:

- 1) We propose the problem of multi-view learning for multivariate performance measures. Although there are plenty of multivariate performance measures optimization methods, they are limited to single view data. There are also lots of multi-view learning methods, however, none of them are proposed to optimize multivariate performance measures.
- 2) We also propose a novel method to solve this problem. We proposed to learn a linear discriminant functions for each view, and combine them to construct a overall multivariate mapping function to predict the class label tuple of a tuple of data points. To learn the linear discriminant functions parameters of different views, we formulate a constrained minimization problem. In this problem, we proposed to minimize the complexity of each linear discriminant functions parameter by minimizing its squared ℓ_2 norm, encourage the consistence among different views by minimizing the squared ℓ_2 norm distance

among different linear discriminant functions response of different views, and also minimize the losses based on a specific multivariate performance measures. The minimization problem is optimized by a cutting-plane method in an alternative algorithm [8]. We maintain a active set of constrains, and update it in each iteration by add a most violated class label tuple. We also optimize the multivariate mapping function parameters one by one, and the optimization of a multivariate mapping function parameter can be solve as a quadratic program problem.

The rest parts of this paper are organized as follows: in section II, we introduce the proposed method, and in section III, we conclude this paper.

II. PROPOSED METHOD

A. Problem formulation

We assume we have a training data set of n data points, and each data point has m views. The training set is denoted as $\{(\mathbf{x}_i^j|_{j=1}^m, y_i)\}_{i=1}^n$, where $\mathbf{x}_i^j \in \mathbb{R}^{d_j}$ is the d_j -dimensional feature vector of the j -th view of the i -th data point, and $y_i \in \{+1, -1\}$ is the binary class label of the i -th data point. We propose to learn a multivariate mapping function \bar{h} to map a tuple of n data points of m views, $\bar{\mathbf{x}} = (\mathbf{x}_1^j|_{j=1}^m, \dots, \mathbf{x}_n^j|_{j=1}^m)$, to tuple of n class labels, $\bar{y} = (y_1, \dots, y_n)$. To implement this multivariate mapping function, we use a linear discriminant function $f_j(\bar{\mathbf{x}}^j, \bar{y}')$ to predict the response of a view the the data tuple, $\bar{\mathbf{x}}^j = (\mathbf{x}_1^j, \dots, \mathbf{x}_n^j)$, against a candidate class label tuple $\bar{y}' = (y'_1, \dots, y'_n)$,

$$f_j(\bar{\mathbf{x}}^j, \bar{y}') = \sum_{i=1}^n \mathbf{w}_j^\top \Psi(\bar{\mathbf{x}}^j, \bar{y}') \quad (1)$$

where $\mathbf{w}_j \in \mathbb{R}^{d_j}$ is a parameter vector for the linear discriminant function of the j -th view, and $\Psi(\bar{\mathbf{x}}^j, \bar{y})$ is a function which can returns a vector to describe the match between $\bar{\mathbf{x}}^j$ and \bar{y}' , defined as follows,

$$\Psi(\bar{\mathbf{x}}^j, \bar{y}') = \sum_{i=1}^n y'_i \mathbf{x}_i^j. \quad (2)$$

Based on the linear discriminant functions of different views, we construct the multivariate mapping function as follows,

$$\bar{h}(\bar{\mathbf{x}}) = \arg \max_{\bar{y}' \in \mathcal{Y}} \left\{ \sum_{j=1}^m f_j(\bar{\mathbf{x}}^j, \bar{y}') = \sum_{j=1}^m \mathbf{w}_j^\top \Psi(\bar{\mathbf{x}}^j, \bar{y}') \right\} \quad (3)$$

where $\mathcal{Y} = \{+1, -1\}^n$ is a set of all admissible label vectors.

We propose to optimize a complex multivariate performance measure, Δ , by learning the parameter vectors $\mathbf{w}_j|_{j=1}^m$ for the m views. To this end, we consider the following three problems,

- 1) **Reducing the complexity of each linear discriminative function:** To prevent over-fitting, we propose to reduce the complexity of the linear discriminative function of each view by minimizing the squared ℓ_2 norm of its parameter,

$$\min_{\mathbf{w}_j|_{j=1}^m} \frac{1}{2} \sum_{j=1}^m \|\mathbf{w}_j\|_2^2. \quad (4)$$

- 2) **Encouraging consistences among different views:** To encourage the consistences of different views, we proposed to minimize the squared ℓ_2 norm distances of responses of any two linear discriminative functions over one single data point,

$$\min_{\mathbf{w}_j|_{j=1}^m} \frac{1}{2} \sum_{j,j':j'<j} \left(\sum_{i=1}^n \|\mathbf{w}_j^\top \mathbf{x}_i^j - \mathbf{w}_{j'}^\top \mathbf{x}_i^{j'}\|_2^2 \right). \quad (5)$$

- 3) **Optimizing multivariate performance measures:** To optimize a specific multivariate performance measure, we propose to minimize the upper boundary of a loss function Δ based on this multivariate performance measure,

$$\begin{aligned} \min_{\mathbf{w}_j|_{j=1}^m, \xi} \quad & \xi \\ \text{s.t. } \forall \bar{y}' \in \mathcal{Y}/\bar{y} : \quad & \sum_{j=1}^m \mathbf{w}_j^\top [\Psi(\bar{\mathbf{x}}^j, \bar{y}) - \Psi(\bar{\mathbf{x}}^j, \bar{y}')] \geq \Delta(\bar{y}', \bar{y}) - \xi, \end{aligned} \quad (6)$$

where ξ is a slack variable of the upper boundary of the loss function.

The overall optimization function are obtained by combining all the problems above,

$$\begin{aligned} \min_{\mathbf{w}_j|_{j=1}^m, \xi} \quad & \left\{ \frac{1}{2} \sum_{j=1}^m \|\mathbf{w}_j\|_2^2 + C_1 \xi \right. \\ & \left. + \frac{C_2}{2} \sum_{j,j':j'<j} \left(\sum_{i=1}^n \|\mathbf{w}_j^\top \mathbf{x}_i^j - \mathbf{w}_{j'}^\top \mathbf{x}_i^{j'}\|_2^2 \right) \right\} \quad (7) \end{aligned}$$

$$\text{s.t. } \forall \bar{y}' \in \mathcal{Y}/\bar{y} :$$

$$\sum_{j=1}^m \mathbf{w}_j^\top [\Psi(\bar{\mathbf{x}}^j, \bar{y}) - \Psi(\bar{\mathbf{x}}^j, \bar{y}')] \geq \Delta(\bar{y}', \bar{y}) - \xi.$$

In the objective of this problem, the first term is to reduce the complexity of the parameter vector of each view, and second term is a slack variable to represent the upper boundary of the multivariate performance measure, and the third term is to encourage the consistences of the responses of different views. C_1 and C_2 are tradeoff parameters.

B. Problem optimization

To solve this problem, we employ the cutting-plane algorithm. In an iterative algorithm, we update the parameter vectors $\mathbf{w}_j|_{j=1}^m$ and an active set of constrains \mathcal{W} alternately.

1) *Updating $\mathbf{w}_j|_{j=1}^m$* : When we have a given active set of constrain $\mathcal{W} \subseteq \mathcal{Y}/\bar{y}$, and optimize the parameter vectors, we optimize \mathbf{w}_j one by one, i.e., when \mathbf{w}_j is optimized, $\mathbf{w}_{j'}|_{j' \neq j}$ is fixed. The following optimization is obtained in this case,

$$\begin{aligned} \min_{\mathbf{w}_j, \xi} & \left\{ \frac{1}{2} \|\mathbf{w}_j\|_2^2 + C_1 \xi \right. \\ & \left. + \frac{C_2}{2} \sum_{j': j' < j} \left(\sum_{i=1}^n \|\mathbf{w}_j^\top \mathbf{x}_i^j - \mathbf{w}_{j'}^\top \mathbf{x}_i^{j'}\|_2^2 \right) \right\} \\ \text{s.t. } & \forall \bar{y}' \in \mathcal{W}: \\ & \mathbf{w}_j^\top [\Psi(\bar{x}^j, \bar{y}) - \Psi(\bar{x}^j, \bar{y}')] \\ & \geq \Delta(\bar{y}', \bar{y}) - \sum_{j': j' \neq j} \mathbf{w}_{j'}^\top [\Psi(\bar{x}^{j'}, \bar{y}) - \Psi(\bar{x}^{j'}, \bar{y}')] - \xi. \end{aligned} \quad (8)$$

The Lagrange function of this problem is

$$\begin{aligned} \mathcal{L}(\mathbf{w}_j, \xi, \alpha_{\bar{y}'} |_{\bar{y}': \bar{y}' \in \mathcal{W}}) \\ = & \frac{1}{2} \|\mathbf{w}_j\|_2^2 + C_1 \xi \\ & + \frac{C_2}{2} \sum_{j': j' < j} \left(\sum_{i=1}^n \|\mathbf{w}_j^\top \mathbf{x}_i^j - \mathbf{w}_{j'}^\top \mathbf{x}_i^{j'}\|_2^2 \right) \\ & - \sum_{\bar{y}': \bar{y}' \in \mathcal{W}} \alpha_{\bar{y}'} \left(\mathbf{w}_j^\top [\Psi(\bar{x}^j, \bar{y}) - \Psi(\bar{x}^j, \bar{y}')] - \Delta(\bar{y}', \bar{y}) \right. \\ & \left. + \sum_{j': j' \neq j} \mathbf{w}_{j'}^\top [\Psi(\bar{x}^{j'}, \bar{y}) - \Psi(\bar{x}^{j'}, \bar{y}')] + \xi \right) \\ = & \frac{1}{2} \mathbf{w}_j^\top \Omega \mathbf{w}_j + C_1 \xi - \mathbf{w}_j^\top \beta + \frac{C_2}{2} \sum_{j': j' < j} \sum_{i=1}^n \mathbf{w}_j^\top \mathbf{x}_i^j \mathbf{x}_i^{j'}{}^\top \mathbf{w}_{j'} \\ & - \sum_{\bar{y}': \bar{y}' \in \mathcal{W}} \alpha_{\bar{y}'} (\mathbf{w}_j^\top \gamma_{\bar{y}'} - \delta_{\bar{y}'} + \xi) \end{aligned} \quad (9)$$

where $\alpha_{\bar{y}'} \geq 0$ is the Lagrange multiplier for the constrain $\mathbf{w}_j^\top [\Psi(\bar{x}^j, \bar{y}) - \Psi(\bar{x}^j, \bar{y}')] \geq \Delta(\bar{y}', \bar{y}) - \sum_{j': j' \neq j} \mathbf{w}_{j'}^\top [\Psi(\bar{x}^{j'}, \bar{y}) - \Psi(\bar{x}^{j'}, \bar{y}')] - \xi$, and

$$\begin{aligned} \Omega &= \left(I + C_2 \sum_{j': j' < j} \sum_{i=1}^n \mathbf{x}_i^j \mathbf{x}_i^{j'}{}^\top \right), \\ \beta &= C_2 \sum_{j': j' < j} \sum_{i=1}^n \mathbf{x}_i^j \mathbf{x}_i^{j'}{}^\top \mathbf{w}_{j'}, \\ \gamma_{\bar{y}'} &= [\Psi(\bar{x}^j, \bar{y}) - \Psi(\bar{x}^j, \bar{y}')] \\ \delta_{\bar{y}'} &= \left(\Delta(\bar{y}', \bar{y}) - \sum_{j': j' \neq j} \mathbf{w}_{j'}^\top [\Psi(\bar{x}^{j'}, \bar{y}) - \Psi(\bar{x}^{j'}, \bar{y}')] \right). \end{aligned} \quad (10)$$

This optimization problem can be transferred to its dual form,

$$\begin{aligned} \max_{\alpha_{\bar{y}'} |_{\bar{y}': \bar{y}' \in \mathcal{W}}} \min_{\mathbf{w}_j, \xi} \mathcal{L}(\mathbf{w}_j, \xi, \alpha_{\bar{y}'} |_{\bar{y}': \bar{y}' \in \mathcal{W}}) \\ \text{s.t. } \forall \bar{y}', \bar{y}' \in \mathcal{W}: \alpha_{\bar{y}'} \geq 0. \end{aligned} \quad (11)$$

To obtain the optimal points of \mathbf{w}_j and ξ , we set the gradients of Lagrangian function with respect to \mathbf{w}_j and ξ to zeros, and we have,

$$\begin{aligned} \frac{\partial \mathcal{L}}{\partial \mathbf{w}_j} &= \Omega \mathbf{w}_j - \beta - \sum_{\bar{y}': \bar{y}' \in \mathcal{W}} \alpha_{\bar{y}'} \gamma_{\bar{y}'} = 0 \\ \Rightarrow \mathbf{w}_j &= \Omega^{-1} \left(\beta + \sum_{\bar{y}': \bar{y}' \in \mathcal{W}} \alpha_{\bar{y}'} \gamma_{\bar{y}'} \right), \\ \frac{\partial \mathcal{L}}{\partial \xi} &= C_1 - \sum_{\bar{y}': \bar{y}' \in \mathcal{W}} \alpha_{\bar{y}'} = 0 \\ \Rightarrow \sum_{\bar{y}': \bar{y}' \in \mathcal{W}} \alpha_{\bar{y}'} &= C_1. \end{aligned} \quad (12)$$

By substituting these results back to (11), we obtain the dual problem as

$$\begin{aligned} \max_{\alpha_{\bar{y}'} |_{\bar{y}': \bar{y}' \in \mathcal{W}}} & \left\{ -\frac{1}{2} \left(\beta + \sum_{\bar{y}': \bar{y}' \in \mathcal{W}} \alpha_{\bar{y}'} \gamma_{\bar{y}'} \right)^\top \Omega^{-1} \left(\beta + \sum_{\bar{y}': \bar{y}' \in \mathcal{W}} \alpha_{\bar{y}'} \gamma_{\bar{y}'} \right) \right. \\ & \left. - \sum_{\bar{y}': \bar{y}' \in \mathcal{W}} \alpha_{\bar{y}'} \delta_{\bar{y}'} \right\} \\ \text{s.t. } & \forall \bar{y}', \bar{y}' \in \mathcal{W}: \alpha_{\bar{y}'} \geq 0, \text{ and } \sum_{\bar{y}': \bar{y}' \in \mathcal{W}} \alpha_{\bar{y}'} = C_1. \end{aligned} \quad (13)$$

We can solve this problem as a quadratic program problem.

2) *Updating \mathcal{W}* : To update \mathcal{W} , we first find the most violated \bar{y}' and then add it the \mathcal{W} . \bar{y}' is obtained as

$$\begin{aligned} \bar{y}' = \arg \max_{\bar{y}'' \in \mathcal{Y}/\bar{y}} & \left\{ \Delta(\bar{y}'', \bar{y}) + \sum_{j=1}^m \mathbf{w}_j^\top \Psi(\bar{x}^j, \bar{y}'') \right. \\ & \left. - \sum_{j=1}^m \mathbf{w}_j^\top \Psi(\bar{x}^j, \bar{y}) \right\}. \end{aligned} \quad (14)$$

3) *Alternative algorithm*: The proposed iterative algorithm is given in Algorithm 1.

III. CONCLUSION

In this paper, we propose the problem of multi-view learning for multivariate performance measures, and a novel algorithm to solve it. Multivariate performance measures optimization can obtain a predictor which directly optimize a desired performance measure. However, the existing methods are limited to single view data, and cannot utilize multi-view data. We propose to learn from multi-view data to optimize the multivariate performance measures, and it has potential in real-world applications. The proposed method can be implemented easily due to its employment of cutting plane method and quadratic program.

Algorithm 1 Iterative multi-view learning algorithm for multivariate performance measures.

Input: Multi-view training data set $\{(\mathbf{x}_i^j|_{j=1}^m, y_i)\}_{i=1}^n$;
Input: Tradeoff parameters C_1 and C_2 ;
Input: The maximum iteration number T .
Initialize linear discriminative function parameter \mathbf{w}_j^0 for each view;
 $\mathcal{W} \leftarrow \emptyset$;
Calculate Ω as in (10);
for $t = 1, \dots, T$ **do**
 Find the most violated class label tuple \bar{y}^t as in (14) by fixing $\mathbf{w}_j^{t-1}|_{j=1}^m$;
 Add \bar{y}^t to the constrain active set \mathcal{W} , $\mathcal{W} \leftarrow \mathcal{W} \cup \{\bar{y}^t\}$;
 for $j = 1, \dots, m$ **do**
 Update \mathbf{w}_j^t by solving (13) by fixing $\mathbf{w}_{j'}^{t-1}|_{j' \neq j}$ and \mathcal{W} ;
 end for
end for
Output: $\mathbf{w}_j^T|_{j=1}^m$.

REFERENCES

- [1] T. Joachims, “A support vector method for multivariate performance measures,” in *Proceedings of the 22nd international conference on Machine learning*. ACM, 2005, pp. 377–384.
- [2] N. Li, I. W. Tsang, and Z.-H. Zhou, “Efficient optimization of performance measures by classifier adaptation,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 35, no. 6, pp. 1370–1382, 2013.
- [3] Q. Mao and I.-H. Tsang, “A feature selection method for multivariate performance measures,” *Pattern Analysis and Machine Intelligence, IEEE Transactions on*, vol. 35, no. 9, pp. 2051–2063, 2013.
- [4] V. Sindhwani and D. S. Rosenberg, “An rkhs for multi-view learning and manifold co-regularization,” in *Proceedings of the 25th international conference on Machine learning*. ACM, 2008, pp. 976–983.
- [5] S. Z. Li, L. Zhu, Z. Zhang, A. Blake, H. Zhang, and H. Shum, “Statistical learning of multi-view face detection,” in *Computer Vision/ECCV 2002*. Springer, 2002, pp. 67–81.
- [6] C. Christoudias, R. Urtasun, and T. Darrell, “Multi-view learning in the presence of view disagreement,” *arXiv preprint arXiv:1206.3242*, 2012.
- [7] W. Li, L. Duan, I. W.-H. Tsang, and D. Xu, “Co-labeling: A new multi-view learning approach for ambiguous problems.” in *ICDM*, 2012, pp. 419–428.
- [8] J. E. Kelley, Jr, “The cutting-plane method for solving convex programs,” *Journal of the Society for Industrial & Applied Mathematics*, vol. 8, no. 4, pp. 703–712, 1960.